
Intelligent tool for visual data analysis



Bachelor's Thesis

Jorge Osés Grijalba

Double Major in Mathematics and Computer Science

Computer Science Faculty

Complutense University of Madrid

May 2019

Intelligent tool for visual data analysis

Jorge Osés Grijalba

Directed by the Doctors

Belén Díaz Agudo, Juan Antonio Recio García

**Double Major in Mathematics and Computer Science
Computer Science Faculty
Complutense University of Madrid**

May 2019

Abstract

*There is light at the end of the tunnel...
hopefully it's not a freight train.*

M. Carey

Index terms— CBR, data visualization, report, IA, artificial intelligence

This document reflects my Bachelor's Thesis corresponding to the Double Degree in Mathematics and Computer Science, developed within the area of intelligent data analytics and 'Case Based Reasoning'. During the progress of the project, the principles applicable in any environment of data processing and the science behind it are explained generally and aimed to be usable in any kind of context by any user provided the right format of data. Nowadays, highly heterogeneous data collection and processing methods are employed in all industries, however the techniques employed to get useful information out of the data usually have a generalistic aim, and the work relevant to the field itself is often done manually. In this work we aim to provide an automated way to analyze information while taking into account information and techniques relevant to the field of the analysis. The objective of this Degree's Final Project is the development of a prototype capable of carrying this analysis while being able to learn based on user input. As a Proof of Concept, we have included several medical domains with each one having developed specific methods and techniques for them. To serve as a base for this analysis, we have also developed a system for storing, loading and analyzing the information of the domain and the information provided by the user. This system will be the backbone of our architecture and enable the Case Based Reasoning analysis to function correctly in very different situations, providing the metrics and functions needed for every case. You can find the code and proofs of concept here.

Contents

Abstract	v
1 Introduction	1
1.1 The Initial Problem	1
1.2 Solution Proposed	2
1.3 Workflow	3
1.4 Structure	4
1.5 Related Work	4
1.6 CBR Design	5
1.7 CBR Process	5
1.8 Retrieve	6
1.9 Reuse	6
1.10 Revise	7
1.11 Retain	7
1.12 Conclusion	7
2 The program structure	9
2.1 The Structure	9
2.2 Overview and Case Presentation	10
2.3 Handling the Information	11
2.4 Analyzing the Information	12
2.5 Presentation and Feedback	13
2.6 Conclusion	14
3 The Program Implementation : Architecture	15
3.1 The Program Module Structure	15
3.2 Non relational databases and the JSON structure	16
3.3 The Objective Storage Module	18
3.4 The Profile Storage Module	18
3.5 The Comparison Metrics	19
3.6 The Analysis Module	20

3.7	The Report Generation Module	21
3.8	The Frontend Module	21
3.9	Information Storage and the CBR	22
3.10	Retrieve	23
3.11	Reuse	23
3.12	Revise	23
3.13	Retain	24
3.14	Conclusion	24
4	Seed Cases and ELO Tournament	27
4.1	Motivation	27
4.2	Seed Cases	28
4.3	Experiment Design	28
4.4	The Elo Rating System	29
4.5	The Tournament	30
4.6	Limitations	30
4.7	Conclusion	31
	Bibliography	33

Chapter 1

Introduction

ABSTRACT: In this study we propose a CBR-based process to retrieve and personalize the visuals from a structured representation of both users and data analysis.

1.1 The Initial Problem

The need for performing analysis on large amounts of data to get very precise and specific information is becoming more and more present everyday in the jobs of data analysts and scientists in every field of work. Large amounts of time are wasted on repetitive tasks such as data wrangling, data transforming and the generation of tailored reports or collections of information with different objectives for diverse profiles with varying degrees of expertise.

These reports are usually formed by a piece of text accompanied by some graphs. It is very common that from these huge amounts of data we want to extract some precise and relevant information to be presented to someone. Our reports have a process behind them that entails the filtering, transformation and selection of the relevant information that will finally be part of the report.

To generate a report we have two questions to answer. First, we must choose the information to present, which is equivalent to choosing what information from the almost unlimited attributes that our data has is relevant to the user. It is clear that this has an *objective* part, in the sense that it is first and foremost a matter of which data is relevant within a certain domain of knowledge and a certain set of metrics, but it also has a *subjective* side, because it's not the same to present a report with medical information to a patient or to present it to a doctor. From this kind of situation the necessity for the *profile* system will arise and we will get both the *doctor* and

patient profiles. The second answer, and perhaps the most *subjective* is how to present it. This decision is related to things like choosing a type of graph, its colors, the font of the text, the words used... and almost an infinite list of *subjective* choices that build upon the previous more objective selection of information to form the final concept of a report of a piece of information.

In this work we propose a system to generate *reports* taking all of these factors into account. In the next section we will flesh out this solution and discuss it further.

1.2 Solution Proposed

From the problem analysis we have concluded that there is no unique formula to generate each report, because it would require the abstraction of very different problems in very different situations and for an almost unlimited variety of users. Furthermore, what if we have to generate a report for a new user? Could this be similar to other reports presented for other users?

Most of these questions can't be answered by a rigid mathematically formulated system, and are best tackled by a mixed system that combines an objective analysis of the information through a variety of metrics, analysis of correlations, distributions and other objective metrics with a subjective approach that takes the final user into account. In this system the knowledge provided by expert knowledge in the field forms the basis for the *objective* approach, while we use an approach based on *experience* of the system to flesh out the final report complemented with the *subjective* side.

Instead, what we propose is a mixed system in which an expert provides an initial input that signals some of the important aspects of the information, and then a pool of experts validates the subjective way in which an information is presented to end up choosing a default report *template* to present to a new user of their *profile*.

To start, we categorize the users or people that will be presented with the report into *profiles*. Then, these groups will provide knowledge of the relevant objective information that they're looking for in the data, like what analysis to perform or what values of certain metrics they'd consider to be relevant. Once this information is fed to the system, it's able to generate reports completing the subjective decisions from semi-random choices from a pool of computer generated graphs, color choices and text based reports. This will provide the basis for a *Case Based Reasoning* system, or *CBR*.

Then our pool of experts proceeds to validate the best report by a voting system based on an ELO tournament, 4.5 which is the case acquisition step of our *CBR* system. The best selected report is then retrieved, reused and presented as default to users of this class.

Each user will also be able to change the result presented to them in

terms of both content (objective) and presentation (subjective), and because the system recognizes individual users it will remember their choices. Users will also have a profile attached to them containing relevant data to the presentation of the information that will allow us to define a metric of sorts between users, the *similarity metric* of our CBR system, to further use this single user customization to influence how information is presented to users of the same class once enough individual inputs have been recorded, possibly substituting the initial ELO based report which will always act as default, thus completing the *learning* step of our CBR system.

1.3 Workflow

The logical structure chosen for the program reflects the need for our tool to be a fully functional agent in and out of itself. We have designed it with a clear divider between a backend capable of storing the information and handling at the lowest possible level, which provides the frontend side with easy methods to get the information it needs, which is then processed taking who is going to look at it into account and then adequately presented to the user.

A cornerstone of the program's functionality is to be able to remember decisions taken by a certain user and to be able to compare new data to old data of its kind.

From these two necessities it is natural to consider some kind of identification system for our datasets, as automated as possible so it needs minimum user input and remains independent of the use case.

For our program, if two datasets contain the exact same set of column names then they are considered to be comparable to each other, and every information stored about this kind of datasets will carry an identifier with the column names.

From here onwards, the term *domain* shall refer to information coming from the same kind of dataset.

When a new dataset doesn't match any previous knowledge, our program automatically creates a new representation for these datasets which is stored along the others. If it detects a matching JSON with knowledge of its domain it loads that instead.

Each representation of a domain stores data such as how many datasets have been loaded and a number of stats for each dataset and its columns depending on its types which will be specified later.

Also, each domain has a number of 'profiles' associated which correspond to *who* this data is associated with. These profiles contain both historical data of the specific owner of the data (in our practical example, the patient data), and who will watch the report generated by this program, that is, the

user of the program.

The information that we're using will be stored in a specific JSON format for each kind that will be specified in chapter 2.

When a dataset is introduced, the program loads the previous information, analyzes it, compares it and generates relevant information to the user. Then it updates the information with both the results of the analysis and user provided information.

This workflow will be the basic use case of the program for every kind of data.

1.4 Structure

A clear module structure is provided so each module does a task in the workflow.

The main modules on the backend side of our application are the Storage module, the CBRStorage module and the Analysis module. For the frontend, the logic structure will be split into the Reporter module and the Presentation module.

Our programming language of choice has been Python, particularly making use of its class to dictionary representation methods which make the work of manipulating the JSON structures much easier than using more rigid languages. A public repository has been created here, and we have used Jupyter Notebooks for the testing and formation of a prototype which has been then moved into standard Python packages.

1.5 Related Work

Natural Language Generation (NLG) produces text in natural language from structured data, the same way we take the structural representations of data and users and generate a report. Our approach relates with research on NLG based on templates (13).

The underlying idea is that texts often follow conventionalized patterns, that can be encapsulated in *schemas*, which are template programs which produce text plans. Schema are derived from a target text corpus, by breaking up these texts into messages, and trying to determine how each message can be computed from the input data. The schema-based approach to NLG approaches problem solving in the very same way as CBR does, reusing previous solutions to be used in future problems, resulting in very much the same cycle of work that we're aiming for.

In (19) we can see a CBR approach reusing explanation reports from

previous transportation incidents. For other related work, (15) identifies core templates to good data journalism practice. In (17) authors describe an approach and a prototype to improve data driven knowledge transfer in presentation tools by applying information visualization concepts. The shortcomings of the current presentation tools are discussed and they also expose narrative visualization, with highly interactive components.

Other uses for CBR have also been used for poetry generation (4) and story plot generation (6). Many of the existing works are presentation oriented (10; 17). CBR has also been successfully used in some help desk applications like the Compaq SMART system(14) and has had a specially successful history in the health sciences(2).

Our approach is more oriented towards visual narrative, i.e., the rendering of data storytelling and visualization of the input data, in the final form of a *report*. Our use case also has the basis on the comparison of new info, represented by the dataset, to old info, represented by both the objective domain information and the subjective information associated to the profile.

1.6 CBR Design

Let's first talk about what we mean when we talk about a CBR system.(16)

CBR, or Case-Based Reasoning, is the process of solving new problems based on the solutions of similar problems we may have encountered before. We can see it as a type of analogy solution making. It draws inspiration from the work of Robert Schank in the 1980s when he developed an early mode for Dynamic Memory(18). It certainly inspired the early CBR systems of CYRUS(9) and IPP(11). All of this eventually resulted in the successful deployment of systems like Clavier(12).

It doesn't need an explicit domain model and so it becomes a task of gathering case histories.

We achieve the reduction of the implementation to essentially identifying significant features that define a case, which is in essence a lot easier than creating a model explicitly.

CBR systems basically learn by acquiring new knowledge as cases, which combined with data handling techniques and big data make maintaining large volumes of information easier.

1.7 CBR Process

Case-based reasoning can be formulated for a program to emulate as the process that follows (1):

1. Retrieve: When facing a new problem, get cases relevant to it from memory. A case is problem, solution, and, optionally, annotations about how the solution was derived.
2. Reuse: Map the solution from the previous case to the target problem.
3. Revise: Having mapped the previous solution to the target situation, test the new solution in the real world (or a simulation) and, if necessary, revise.
4. Retain: After the solution has been successfully adapted to the target problem, store the resulting experience as a new case in memory.

1.8 Retrieve

Conceptually, we need to get the necessary information about past cases of how they were resolved, that is, mainly what problem it was and how it was solved. We start from the idea that our problem is basically analyzing a new dataset. To do this we provide the frame of domains, which ensures us that what we retrieve was relevant condensed information about the problem in the past, and within that information we have the metrics used to analyze datasets like our current one. Another side of how to solve it is represented by the profile information, which tells us how to solve it for the specific user who is using the program.

1.9 Reuse

The knowledge base is obtained primarily from the enumeration of certain past cases or problems. This is built from the fact that experts (humans) are much better at recalling previous experiences and problems than at creating systems of rules.

As new problems are fed to our expert system (containing the knowledge or memory of previous experiences) to which no past problem can match exactly, the system is capable of reasoning from more general similarities to come up with an answer.

This tries to imitate the generalization capability of humans.

To map the solution from the previous cases to the current problem what we do is run an analysis on the current dataset, and then use the metrics contained in the domain information, and then apply the profile information on the analysis to filter the results.

1.10 Revise

When the users are presented with the new report, they have the ability to modify its contents, both graphical and the information they're presented with. This changes are stored in the profile information, so we have access to the new preferences, thus making a better solution from the old one in the spirit of the revise process.

What we're trying to achieve is giving the user as much power as possible, because that is what will make his or her experience better with each time they use the program, which is the whole point behind the implementation of the CBR based system that we've chosen.

1.11 Retain

The system is able to retain new information by being able to make changes to the elements it stores, namely updating the domain objective analysis with the new analysis and updating the subjective profile report changes with the changes introduced by the users.

This is done only once, after the user has made all the changes and exits the program. No interaction is needed on the side of the user to do this, it is done by default so as to make the whole experience as streamlined as possible.

The memory system present in our system grows and changes by each time we present it with a new case. An important aspect of this memory-based process of reasoning is closely related to automatic learning: our system should be able to remember the problems that it has been presented with and to use past information to solve new challenges.

This is intended to complete our modeling of the human behaviour of CBR, and represent the final step of our CBR system.

1.12 Conclusion

In this chapter we've outlined the most important aspects of the program, the related work that has been done in the field and explained the CBR process. We've proposed a workflow to follow, wrapped around a heavy CBR-like program structure that will allow us to work with different kinds of data and different users. In the next chapter we will provide a more specific description of the program functionalities and components, going over design and programming decisions from a high-level point of view.

Chapter 2

The program structure

ABSTRACT: We will provide an analysis of the program structure from an abstract point of view, going from how our cases are presented to how we handle and analyze the information, ending with an explanation of our subjective side and feedback mechanism, and its relation to the conceptual CBR process.

2.1 The Structure

Our program consists of a series of modules, designed with a core of a CBR structure and a series of auxiliary modules and methods to make it easier to function. These modules are designed to interact with each other and provide the necessary tools to work smoothly regardless of the data being used.

In this chapter we take a broad look at the program, trying to explain the design choices behind each aspect of it and to provide a framework of thinking for the next chapter, when we will dive deeper into the technical, low level side of things.

Our program structure needs to be flexible and adaptable, and fit the needs of both providing a basis on which to develop a CBR system and at the same time be quickly handling data, be able to work with very different domains, and be responsive to user input, and scalable to a big data practical solution that can be used by real users for real use cases.

We will have to take into account all of this into every decision, to make sure that we're building a functional system good enough for scalability but at the same time being quick to develop for the proof of concept and the

experiment to test our design choices that will allow us to try our ideas before scaling.

2.2 Overview and Case Presentation

The input for our program is really just an user who desires to analyze a certain dataset, so it can be conceptualized as a tuple containing a user and a dataset, or a *profile* and *domain* with the new information that it represents. This dataset will also belong to a specific *domain*, which will have association with the user *profile* aswell, so our input will be quickly transformed by the program into a domain, *profile*, dataset concept. The concept of extracting the domain information from the dataset provided is our approach to the retrieve from the case base step from a CBR point of view.

From now on on this chapter, we assume we have previous information for both the domain the dataset belongs to and about the type of user that is trying to analyze it. As an example, let's say that we're a doctor looking to analyze how our new patients are faring compared to our past patients in a certain medical test.

Our use case has the basis on the comparison of new info, represented by the dataset, to old info, represented by both the objective domain information and the subjective information associated to the profile. This flow is highly compatible with the CBR approach that we will outline in chapters 5 and 6, because we're essentially solving new problems from old problems. In the doctor's case, they would feed the system new data regarding our new patients' test.

Tying the new information to the old information will be done through a storage system with matching keys for the dataset and domain, being the number of columns in our proof of concept but being easily expandable into a more global system through database indexes or anything acting like a unique identifier. This will be our retrieval step from a CBR point of view, and in the doctor's case this would be done by identifying what kind of tests the data belongs to. As we see in 3.2 this will be done in our proof of concept by identifying the *domain* associated with the dataset fed to the system.

The analysis of the new information will have two parts, the first one will do an isolated analysis of the new information using the techniques provided by the old, and in the second it will compare the results of the first part with the results of a similar analysis stored in the old information. In the case of our medical data, we might be able to note things like the new patients having a generally higher score for some kind of medical test, lower blood pressure, being older than usual or other metrics of the sort, which would be related to *single columns* of the dataset. We would also be able to see things like the new dataset containing only six patients while the datasets usually contain

twenty or more, objective metrics concerning the *whole dataset*, or even a correlation between age and blood pressure, objective metrics concerning *two or more columns*. This will provide us with a new object representing the relevant information about our new problem, or dataset, and then we will filter this object using the subjective information from the profile to present the user with a final report.

The user will then have the means to modify this report, and thus modifying the subjective information about him or her, forming the *provide* CBR step. Let's say for example that as a doctor we've been presented with a *report* containing among other things a graphic for the unusual blood pressure from our patients in blue and also notes the unusual number of patients aswell. Then we'd be able to tell the program that the number of patients is not really relevant and shouldn't be taken into account, and that we'd also like a different type of graph for the blood pressure. We're able to change the report and when we shut down the program it saves the final state of the report.

The program will also use the new analysis to update the old information about the domain before proceeding to store it. This will provide our program with the *reuse* and *retain* functionalities needed from a CBR perspective.

We will expand on how the base cases are formed and how we adapt to new profiles and new domains on further chapters, but for now let's say that a process is in place to ensure that the default report is correct enough, the analysis performs its due processes and the result is at least acceptable and relevant to the final user.

2.3 Handling the Information

First and foremost we need a way to store information and to retrieve it effectively. We will store two different kinds of information: *objective* information concerning the results of mathematical analysis performed on a dataset and *subjective* information on what information to present and how to present it to the user. Since the first is related to a certain domain and applies to all its associated users (or *profiles*), and the second concerns both the domain and the user, it makes sense to store and handle them separately.

Regarding the first kind, we now provide both a way to identify which domain it belongs to and a general description of what it contains, which will be specified later during the technical implementation chapter. It contains statistics and properties from both columns and datasets, different for different types, as well as how they were measured. Saving how they're measured is important to measure new datasets and to be able to compare metrics effectively.

For a basic approach of our *objective* column metrics we have divided

columns into categorical data, that being columns which have discrete non-ordered values, and *numerical* data which contains numbers. For the first we will measure things like most frequent categories and category distribution, while for numerical columns we have a more varied set of tools like studying ranges, numerical distributions, means, medians and other common mathematical metrics.

Both the column specific stats and the dataset stats are subjective to change with each dataset added to the domain, so it is important that we're able to update this information readily and effectively. For starters we have provided our program with a toolset capable of doing an array of the most common analysis in data science.

It's also important to be able to distinguish between *domains*, that is, between datasets with different sets of columns. The solution proposed is to use the columns of each dataset as a unique key that identifies a domain, and their objective information will be stored as part of the same object.

In the case of subjective kind, we combine this with a unique user identifier to make a unique key to identify the information.

For all of this to be possible the program has to be able to handle loading the information, storing it and updating it, and that is achieved having functions dedicated to loading this information into the program and updating it when presented with new datasets and storing it overwriting the previous information, creating a dynamic system capable of learning from users and gathering new insights.

2.4 Analyzing the Information

The *objective* analysis will concern the objective side only, so we don't need to take the user into account for now.

First a dataset is provided as input for our tool. From this dataset, we retrieve the information of the associated domain so we know how to perform the analysis on this new dataset, as it has to be comparable to the domain information.

So, using the metrics and analysis detailed in the domain information, we perform them on the dataset and compare the results with those of the domain through a series of comparison metrics which will be detailed later.

These results are then condensed into an object containing all the objective data from the comparisons.

This output will be the input for the subjective side of our program, as this information will then be filtered and transformed in a way tailored to the user.

The *subjective* analysis is to be performed after the previous step has

been completed.

From the representation of the objective analysis, we need to perform two tasks. To perform these two tasks we first load the information we know about the user, which has been stored separately from the objective information as we have previously stated. One is the filtering of such information, selecting the information relevant to the user. This is achieved by using the information that the user's *profile* contains, which will tell us the information of the analysis that will be relevant. The second task is choosing how to present this information to the user, that is, using the information also stored in the *profile* to infer the appearance of our final *report*, that is, how will we represent the relevant information to this user.

2.5 Presentation and Feedback

Once the user has been presented with the report, he's able to make changes to it using a graphical interface. When the report is saved and the user has exited the program, it initializes a shutting down routine on which the two information databases are updated with the changes made by the user (if they existed) and the new information provided by the dataset.

This is essential to our CBR process as it creates a way to tie the old with the new, and greatly increases user experience as the program gets better with each iteration, as intended with the CBR methodology.

This is the main way our program has of expanding its knowledge, which will then be used to present better reports to the same user and to adapt the information it presents to new users which will be similar to this one.

So we get a double benefit, getting both better results for this user and for all users of the same domain and class each time someone uses our program.

In terms of the presentation, we have two axis on which to make decisions, visual (graph-like) content and textual content.

For the textual side, for the proof of concept we've developed a very simple template system for variables which lacks context for the specific domain and is instead tied to function results, that is, only distinguishes between categorical and numerical values for the final explanation.

On the visual side, for the proof of concept we've added pie plots, histograms and temporal graphs each with an array of different colors to choose from, because of the importance of colors in the representation of information.

It's also important to note that users are able to change the scale, size and place of the graphs in the report, which is very important for the visual impact and the user's quick and easy understanding of the information that is being presented to them.

For the user, being able to change all of these aspects represents a very important interactive side of the tool, giving power to the user as it empowers their experience with each interaction.

2.6 Conclusion

In this chapter we've provided a more detailed view of our program's usage, going over its functional modules from an abstract point of view. We've also provided some specific description of the program functionality needed to implement at a programming level. In the next chapter we will provide a more specific description of the program functionalities and components, going over design and programming decisions from a low-level point of view, and going over how the different CBR steps were taking into account and implemented.

Chapter 3

The Program Implementation : Architecture

ABSTRACT: In this chapter we provide a technical analysis of our tool, examining how it works and what was designed for at a programming level, as well as its module structure, and its relation to the low-level CBR process.

3.1 The Program Module Structure

First we need to provide our program with a clear programming structure that allows it to be customizable enough, as it is very important that we're able to add new functionality related to new kinds of analysis, users or datasets.

To make it as customizable as possible, we've divided it in modules that have a clear functionality, with the objective that we're able to change an aspect of the program by changing just one module and not the whole program.

All the code is openly hosted here, including the early concepts, data needed for a couple of POC, a graphical interface for a tournament structure, all the classes listed and an early server concept for the tournament.

We will have one module which handles the loading and storing of the objective information about domains, one that does the same for user profiles associated to each domain, one that runs an analysis on the dataset and

compares it to the domain, and one that puts everything together to generate a report which will be then presented to the user through the frontend module, which handles interactions with the user.

So our hierarchy goes as follows, the frontend module takes input from the user regarding its profile (input done through a dictionary-like object in the proof of concept) and a path to the dataset that is to be analyzed.

This module then relays the information to the two modules that handle information storage.

The storage module will read the dataset, as a CSV file in the proof of concept, read its columns and search its database for a matching domain information.

In a more advanced version this search would be performed through a lookup on a non relational database like MongoDB, but for the proof of concept it looks in a folder where it stores the domain information objects as JSONs.

When it finds one that matches the column names, it loads the object into a python dictionary and sends it to the analysis module along with the new dataset.

The analysis module then performs an analysis and generates a proto-report with all the objective information about the comparison between the dataset analysis and the domain historical analysis, then passes this along to the report generation module.

In parallel, the module that handles the profile will load it in a way similar to the storage module, then pass it along to the report generation module.

When the report generation module gets all the needed information uses the user-generated profile to "filter" the objective results (as explained later in detail), and extract from them the final info that will then be presented to the user through the frontend module.

The user is able to make changes to this final report (through the console in the proof of concept), and before closing the program

All of this contributes to making the CBR process as streamlined as possible while providing a flexible framework for data handling.

3.2 Non relational databases and the JSON structure

For the proof of concept implementation developed in this thesis, we've used the JSON structure to handle our information, to write it to a file system and to read it later.

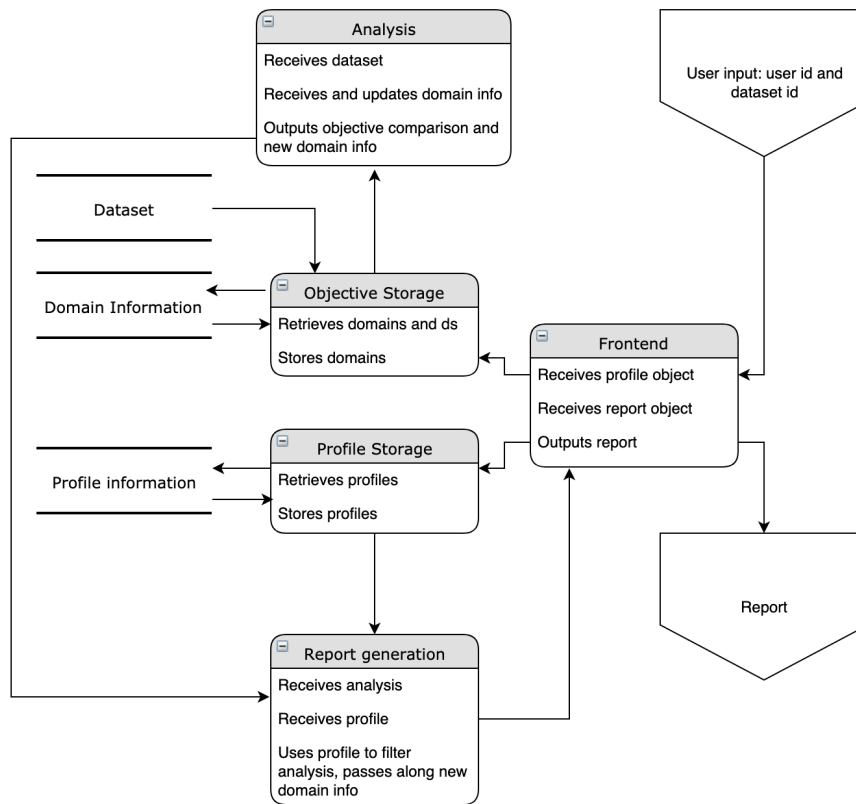


Figure 3.1

JSON, short for JavaScript Object Notation is a file format that uses text readable by humans to write attribute value pairs and serializable values.

It is a very common data format used for general storage of information, and will allow us to do a quick implementation capable of handling enough domains and profiles to test our program, and will allow us to later transfer this information to a non-relational database system such as MongoDB.

Python also allows us to quickly load this kind of object into its native dictionary type, making it easy for our modules to load, manipulate and store the information.

A non relational database is one that does not use the tabular schema of rows and columns found in most typical databases like SQL.

Our need for such a database comes from the fact that not every domain will have the same structure of knowledge for its analysis, with the possibilities of combinations being almost infinite and thus negating any approach to putting it inside a tabular schema.

Inside these databases data may be kept as JSON documents, which is the format that we've chosen.

The proposed database for the final model is MongoDB.

MongoDB is a non relational database program, which is also cross-platform and document oriented. It works with objects very similar to JSON, and is licensed under the Server Side Public License (SSPL).

Aside from this, it's easily managed through Python and has native packages to communicate with it, which would make our life easier when implementing the full program.

All of this provides us with a solution catered to our needs, which can be scalated easily if needed into a production ready state, capable of dealing with large amounts of data and providing a real solution to real users.

Every module listed below corresponds to a Python class, and has a number of methods, functions and dictionaries associated to it to perform its function.

3.3 The Objective Storage Module

This module is capable of reading and updating the information stored about the objective analysis of past domains.

It sits at the lowest point in the hierarchy, as it only does as commanded and is usually handled by the other modules and used as a mere tool to get information, similar to the profile storage module.

Both modules are derived from a more basic Python class, called just Storage, from which they inherit the methods responsible from loading and storing files (or updating and loading from the database, in the scalable version).

This class adds to these functions more advanced tools to deal with the domain information and load it into a dictionary-like object that will then be able to be easily handled by the analysis module, translating information stored in plain text like function names to its equivalent variables and Python dictionaries in the program, using the dictionaries which it has as an attribute.

This process is reversed before writing the information back to disk in a format that can be condensed to JSON documents, mainly turning every function name and other serializable objects back into strings to be stored.

3.4 The Profile Storage Module

Another kind of information is stored. This is the one concerning the human side of things, that is, how to interpret these stats and turn them into something that humans with different levels of familiarity can understand.

To do this, we provide use another storage class that will contain human-relevant data that will modify the objective comparison delivered by the analysis module.

A system of profiles is added to the object itself, inside a "profiles" key. The domain associated is clear as they share the same "attributelist" identification system.

The information contained in each one of these "profiles" serves two different purposes. It provides customizable elements of how the data will be presented to the user, and it keeps an historical record of this user's dataset results (similar to the one in the domain storage) making a historical following of a profile possible.

For each domain, there's a default profile. This provides a way to present the data when no previous knowledge of the profile is available. The automated processes of obtaining this profile and tuning the existing profiles from user feedback will be explained in further chapters.

3.5 The Comparison Metrics

The purpose of these functions is to provide a way to measure the properties of a given dataset or knowledge domain.

We have to note first that not every kind of possible metric has been added to the program, but we've instead added enough to cover the most common types of analysis, mainly detecting distributions, most frequent values, and calculating a series of common metrics over numerical columns.

However, adding a new metric is as simple as providing it with a name, an associated function which fits into the types of one of the metrics listed below, and adding the pair of name,function to a dictionary present in the code.

The rest of the program will behave exactly the same, thus fitting the design principle mentioned before of being able to expand the program quickly and efficiently.

We can categorize them as follows:

First the "measurement" metrics, used to get the information of a single dataset or domain.

- Dataset Metrics : they concern the dataset as a whole, like number of rows with missing values. `dm :: (ds) -> num`
- Single Column Metrics : they concern a certain column, and are based on the type of the column. For numerical columns we will have things like median, averages, deviations, distributions... For categorical columns we'll work with frequencies and things of the sort. `scm :: (col) -> num`

- Multiple Column Metrics : we will be looking for correlations and things of that sort. `scm :: (col,col) -> num` Time based metrics will be defined from this construct.

We will also have "comparison" metrics, used to compare datasets against their domains. These metrics will compare the output of two measurement metrics, both will have to spawn from the same function.

- `compm :: (metric) -> num`

Note that these metrics are not to provide "meaning" or any human-readable input, nor to be inherently comparable between each other outside of a framework of understanding of the domain (metric importance).

A mean to convert these machine cold metrics into human understanding will be provided in further modules. For now, we're not taking humans into account.

3.6 The Analysis Module

The analysis module will receive a dataset, use the Storage module to load its information, then analyze the dataset, which generates a similar object to the domain json, then producing a comparison of both.

The main methods for this module are `getstats` , `getcolumnstats` and `getdatasetstats` .

The `getstats` method is just a wrapper for the other two, calls them both and stores its results inside the Analyzer class.

Both `getcolumnstats` and `getdatasetstats` compute the statistics for the given dataset. If there is previous knowledge of the domain, the stats that appear there are computed for the domain. If not, a standard set of frequencies for categorical values and medians and distributions for numerical values are calculated and used to populate the stats object.

The most important method is `getanalysis` . Once the stats for the datasets have been generated, if there's previous knowledge available the class runs an analysis comparing the metrics of the two, and generating an object with the result. For this to happen, each metric defined for the dataset must have an associated comparison metric .

If there is no previous knowledge then the dataset stats are passed along to the reporter with a field indicating that there was no previous knowledge.

In any case, at this level we've already filtered what is relevant and what is not from the comparison.

3.7 The Report Generation Module

This module stands at the edge between the backend and the frontend. It receives the information from the comparison between the dataset and the domain knowledge and extracts the relevant profile information.

Once this is done, it uses both to generate a report with all the information from both sides. The user-relevant info will modify what is shown and how that is shown, changing the graphical elements according to the user so the frontend modules are able to be logic-free.

It is able to directly modify the profile information by the proxy methods `modify` and `savehumaninfo`. Its main method, `generate`, will create and populate an attribute within itself called `report`.

This method is called when the class itself is generated but the report can be modified as any Python attribute if needed.

At this point, we have an object that represents the dataset compared to the historical data and data about the profile associated with the user. This information, however, is in the form of a JSON object and is not really human-readable. The job of the frontend modules is to take this information and turn it into something easy to understand.

3.8 The Frontend Module

The purpose of this module is to handle the user program interaction. It sits atop of the program hierarchy, being able to use all the other modules as it sees fit.

Its role begins at the very beginning of the program life cycle and ends at the same time the user decides to exit.

Its functionality is based on the principles of minimalistic design and the user being able to interact with every element of the presented report.

For our proof of concept, a more simple design has been put in place. Instead of clicking on the elements, the user is able to interact with the program through the console.

However, all the functionality is present, as the users can input the information the program requests to generate the profile, and then change the report presented to them using commands too.

When the program is closed the frontend module passes the modified report to the information storage modules, with both the updated objective information and the updated subjective information.

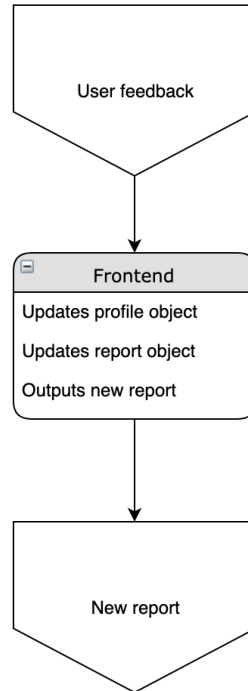


Figure 3.2

3.9 Information Storage and the CBR

Because the information we're storing is still low and the metrics are different for each domain the only way to consistently store the information is through a non-relational database.

To provide the information to the experiment and to run this as a proof of concept we can simplify it by storing information in the form of JSON structures, which will be stored as files in a directory accessible by our program.

The relevant information for the use case at hand will be retrieved by our program, stored as a run time variable, modified when needed and then stored back to the JSON overwriting the previous structure. If we had a non relational database we could simply update the database instead.

3.10 Retrieve

If we're doing a full non-relational database implementation it's probably better to develop a domainID, because that would allow us to have different domains with the same column names.

For the proof of concept though that is not an issue because we've made sure this hasn't happened with the datasets used for the testing.

Retrieving this information is as easy as reading the JSON file and loading into a dictionary type in Python. We know which information to retrieve because each json contains an ID variable with the column list of the datasets from the domain.

We do this with a file buffer the standard way. Once it has been loaded the information is handled by two different classes.

3.11 Reuse

The objective information about the domain is loaded and handled by the class that handles the objective information, which will then be passed to the

The profile or subjective information about the domain and user is handled by the module that handles the profiles, which will then pass it along to the module that generates the report.

This allows us to reuse the previous experiences in the process of generating the report module.

3.12 Revise

There are two ways on which we can revise the information. If the users from the start say what they want then the system associates what they want as if changes have been made to the report that was going to be presented to them.

In any case, when the users are provided with the report they can provide feedback through a visual interface by selecting which information is shown and shouldn't be shown, what information should be shown but is not the report, and change the colors/fonts of the report.

This allows us a thorough revise step, in which the user has full agency over the results generated by the program and is able to change all the choices that our system has made.

This is incredibly important because it makes our system learn a lot from

each interaction.

3.13 Retain

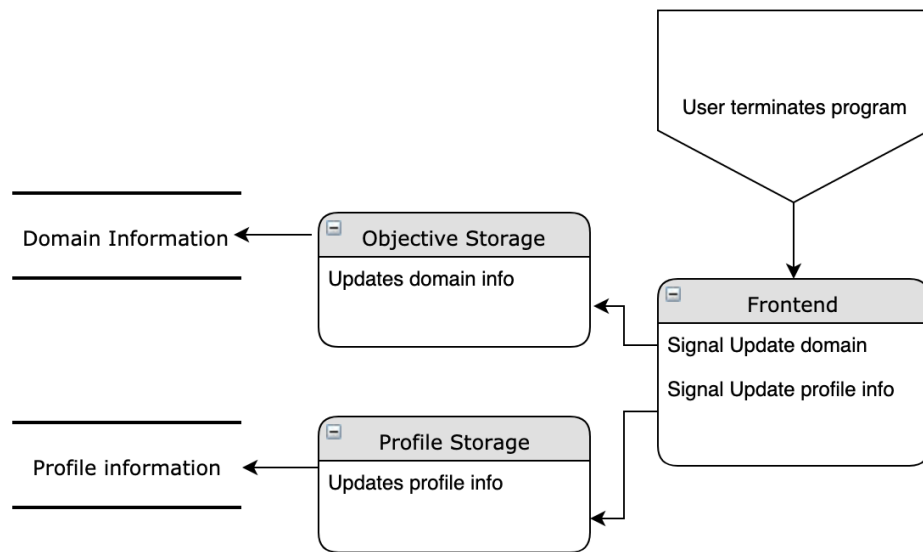


Figure 3.3

Once these changes have been made the objective information is updated and written back to disk by the same module that retrieved it through an store method.

The same method is present in the module responsible for handling the profile information, which will update the choices the user has made.

This is a crucial step because it is what allows us to retrieve these changes consistently if this user or a similar one wants to use the system in the future.

In the end what we've achieved is a system that grows better with each interaction, providing the users with better and easier use experiences every time they use it.

3.14 Conclusion

In this chapter we've gone over the detailed implementation of our program, including the specific modules programmed in Python code and their relations to the flexible CBR structure we're trying to achieve. At this point we

have a clear picture of what our program does and how it looks like from a code perspective, and we've specified the different techniques and metrics used to analyze common data types. However, there is one important question that still doesn't have an answer. How do we generate the first seed cases from which we build our retrieve step? Is there a better way than to randomly assemble them? Can we rely on expert insight for everything? We will answer these questions in the next chapter.

Chapter 4

Seed Cases and ELO

Tournament

ABSTRACT: In this chapter we outline and explain the our techniques used to ensure new users and domains have a valid starting point from which users can productively use the tool.

4.1 Motivation

Our objective for this chapter will be to provide details on the process chosen to develop a solid knowledge basis on which we can build a use case for our program.

What we mean by this basis is the knowledge of which metrics to use to provide a clear picture of the datasets belonging to the domain, as well as which metrics are relevant to a profile and how should they be shown, as with a graph, through a text report, which colors, etc.

We will be using the input of an expert to determine the metrics to use, and will then generate a seed profile based on an ELO tournament with a pool of experts. This will be then affected by user input.

4.2 Seed Cases

The process to establish a frame of knowledge for a certain domain will be initiated by an expert who will provide its associated class, or basic categorization of users, and a brief computer-understood description of what he's interested in.

This knowledge will be used to run a first analysis and generate an array of reports through a semi randomization process, which will be then pooled together as participants of a tournament, and given an initial ELO of 1000.

The final winner of this tournament will be taken as the "seed" for the classification the users belonged to, so if a new user enters the program and belongs to this class, it will be compared to people of this class through a metric and given the report of the person that is closer to them.

The attributes of the user that form the metric are different for each class, and can be things like gender, age, medical specialty, etc.

4.3 Experiment Design

To test this approach to seed generation we have created an

To design the experiment we have used anonymized historical grades data from UCM's Computer Science.

Our dataset contains data about anonymized global degree grades from the nineties, containing the year of graduation and the gender among other variables.

This will be our domain.

This dataset is to be viewed from three different perspectives :

1. Students wishing to know how their grade stands among their peers,
2. Teachers who want to know how the year they've taught has fared compared to the others,
3. The figure of a gender delegate who wishes to know if the grade distributions are different when broken down by gender.

Our objective first is to generate successful seed cases for each class which will then be used as base report generation techniques for each category of user within this domain.

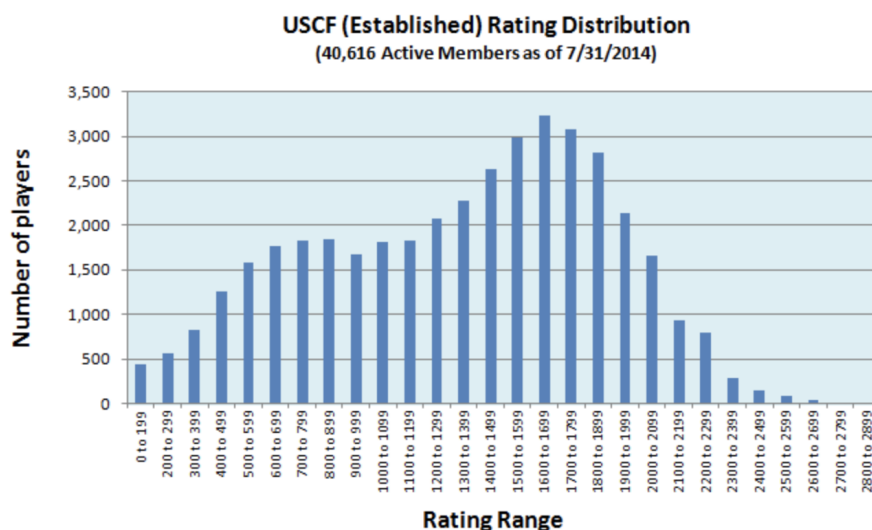


Figure 4.1: ELO for US Chess Federation Ratings (3)

4.4 The Elo Rating System

The Elo rating system is a method for estimating the skill levels of players relative to each other used in zero sum games, most famously chess. It was introduced by Arpad Elo, an active master-level chess player and member of the United States Chess Federation, when looking for a fairer system to replace the numerical rating used before 1960.⁽⁵⁾ It has been adapted to numerous other ratings for specific use cases and situations.⁽⁷⁾ The difference between the scores of two different players serves as a good predictor of the match results.

In our example, the "players" will be the reports generated by our seed generation system, and the winner of the match will be the report chosen by the expert users through the click of a button.

A reports's Elo rating is represented by a number, the higher the better.

We've made our initial rating for every report before initializing the tournament as 1000.

The rating goes up or down based on the result of different matches between reports.

After every game, the winner takes rating away from the loser and onto himself. The difference between the ratings of the winner and loser makes the result vary in quantity, so if a player with a low score wins against a player with a high score the loss is greater for the loser and the gains greater for the winner than in the other situation where the higher rated player wins.

The distance between ratings is also taken into account in this situation, the greater the difference between ratings the bigger the actual gains or loses will be.

This system is thus self correcting and expected to provide a good framework for comparing our reports.

4.5 The Tournament

Our tournament is played in a Round Robin style to ensure that all reports are compared against each other.

We also have to note that the reports are distinguished between one another in different degrees, that is, two can only be differenced by the colors while the difference between two others in a match might entail content.

All in all, we're trying to make very little difference in treatment between them before presenting them to the users, because we assume that our generating system doesn't know the importance of the different aspects of the reports and might take away important information if we make it have more filtering capabilities before the tournament starts.

Elo systems tend to create distributions such as the one in Figure 4.1, giving us a clear winner under normal conditions.

Our first sample for the tournament experiment will be of 12 images, thus creating enough matches for our ELO system to be effective while being small enough not to bother our user pool of experts who will essentially decide the winners of every match

Other proposed alternatives for case generation include completely random seeds that include a large collection of metrics, with the drawback that the first users will have to undergo great customization processes.

4.6 Limitations

With this metric system we lack a way to infer some *broadly* similar information representation characteristics from similar *profiles* in different domains. Although this idea will not be developed further due to time constraints, it would require a way to *categorize* profiles from different domains into the similar categories according to some *inter-domain metric* that would need to be defined. This would allow us to incorporate domains similar to old ones with ease, and would allow us to cover one clear drawback of our proposed system, the fact that when introducing a new domain very similar to an already established one we have no more ground to work on that when introducing a completely different domain, and it shouldn't be the case.

There has been work done recently that adds a statistical framework to the CBR concept and provides a formal case based inference as a probabilistic inference, This way we can generate case based predictions and add a level of confidence to them. (8) This idea would be useful as another approach to relate different domains and see if we can reuse knowledge for similar ones.

4.7 Conclusion

Through this chapter we've covered a first way to answer the questions of how to set the seed cases and obtain a good basis of rules on which to build the case retrievals and similarity metrics of our CBR system, as well as greatly improving its reusability. The main work done through this thesis has been to conceptually develop and also implement a CBR based system capable of handling different types of information, relating them to past information of the same kind, analyzing its differences according to some metrics relevant to its information, and finally using its user knowledge to filter out this information and present it in a relevant way.

Bibliography

- [1] A. Aamodt and E. Plaza. Case-based reasoning: Foundational issues, methodological variations, and system approaches. *AI Commun. IOS Press*, 7(1):39–59, 2016.
- [2] S. Begum, M. U. Ahmed, P. Funk, N. Xiong, and M. Folke. Case-based reasoning systems in the health sciences: A survey of recent trends and developments. *Trans. Sys. Man Cyber Part C*, 41(4):421–434, July 2011.
- [3] M. E. C. Elo distributions for the uscf. 2013.
- [4] B. Díaz-Agudo, P. Gervás, and P. A. González-Calero. Poetry generation in colibri. In *European Conference on Case-Based Reasoning*, pages 73–87. Springer, 2002.
- [5] A. Elo and S. Sloan. *The Rating of Chess Players, Past and Present*. Ishi Press International, 2008.
- [6] P. Gervás, B. Díaz-Agudo, F. Peinado, and R. Hervás. Story plot generation based on CBR. *Knowl.-Based Syst.*, 18(4-5):235–242, 2005.
- [7] M. Glickman. Mathematics of the elo formulation. 2001.
- [8] E. Hüllermeier. *Case-Based Approximate Reasoning (Theory and Decision Library B)*. Springer-Verlag, Berlin, Heidelberg, 2007.
- [9] J. Kolodner. Reconstructive memory: A computer model. *Cognitive Science*, 7:281–328, 1983.
- [10] R. Kosara. Presentation-oriented visualization techniques. *IEEE Computer Graphics and Applications*, 36(1):80–85, 2016.
- [11] M. Lebowitz. Memory-based parsing. *Artificial Intelligence*, 21(4):363–404, 1983.
- [12] W. H. Mark. Case-based reasoning for autoclave management. 1989.

-
- [13] K. R. McKeown. Discourse strategies for generating natural-language text. *Artificial Intelligence* 27(1):1-41., 1985.
 - [14] T. Nguyen, M. Czerwinski, and D. Lee. Compaq quicksource: Providing the consumer with the power of AI. *AI Magazine*, 14(3):50–60, 1993.
 - [15] A. Ojo and B. Heravi. Patterns in award winning data storytelling. *Digital Journalism*, pages 1–26, 2017.
 - [16] M. M. Richter and R. O. Weber. *Case-Based Reasoning: A Textbook*. Springer Publishing Company, Incorporated, 2013.
 - [17] R. Roels, Y. Baeten, and B. Signer. Interactive and narrative data visualisation for presentation-based knowledge transfer. In G. Costagliola, J. Uhomoibhi, S. Zvacek, and B. M. McLaren, editors, *Computers Supported Education*, pages 237–258, Cham, 2017. Springer International Publishing.
 - [18] R. C. Schank. *Dynamic memory - a theory of reminding and learning in computers and people*. Cambridge University Press, 1983.
 - [19] G. Sizov, P. Öztürk, and E. Marsi. Let me explain: Adaptation of explanations extracted from incident reports. *AI Communications*, 30(3-4):267–280, 2017.